arm

# Debugging and Profiling with Arm Forge

ATPESC 2021

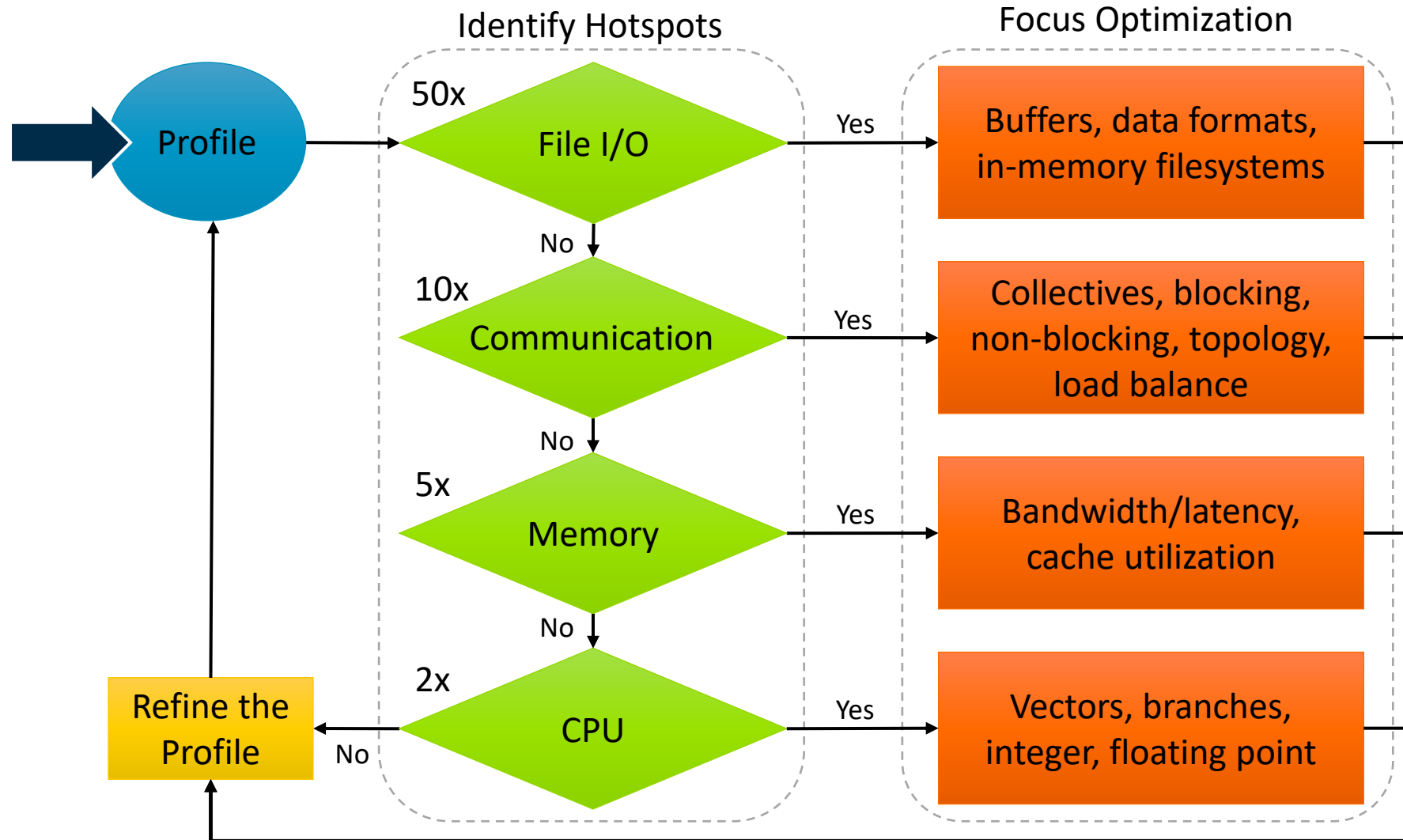Ryan Hulguin
August 11, 2021

# Agenda

- General Debugging and Profiling Advice
- Arm Software for Debugging and Profiling
- Debugging with DDT
- Profiling with MAP
- Theta Specific Settings
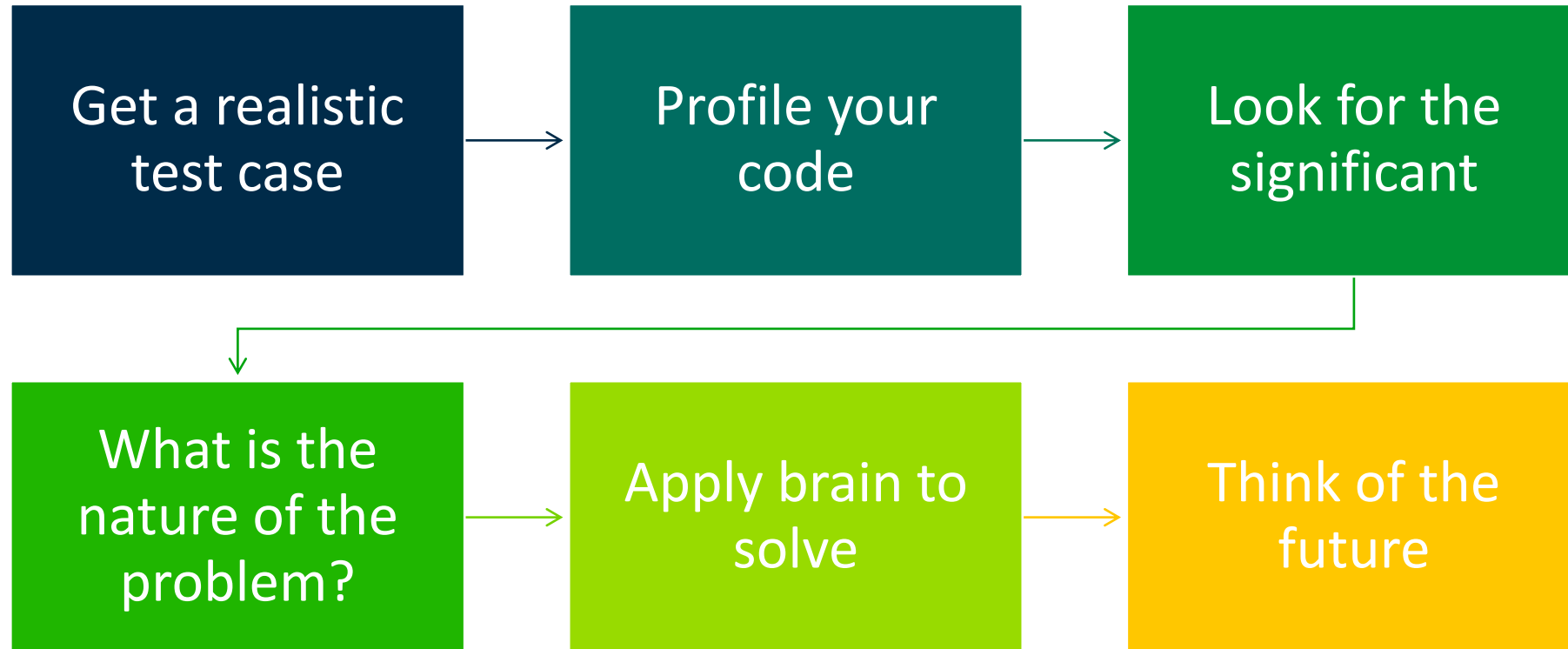
arm

# Debugging

- Transforming a broken program to a working one

- How? TRAFFIC!

  - **T**rack the problem

  - **R**eproduce

  - **A**utomate - (and simplify) the test case

  - **F**ind origins – where could the "infection" be from?

  - **F**ocus – examine the origins

  - **I**solate – narrow down the origins

  - **C**orrect – fix and verify the test case is successful

# Profiling

Profiling is central to understanding and improving application performance.



Identify Hotspots

Focus Optimization

Profile

50x File I/O — Yes → Buffers, data formats, in-memory filesystems

No

10x Communication — Yes → Collectives, blocking, non-blocking, topology, load balance

No

5x Memory — Yes → Bandwidth/latency, cache utilization

No

2x CPU — Yes → Vectors, branches, integer, floating point

No → Refine the Profile

arm

# Performance Improvement Workflow

Get a realistic test case → Profile your code → Look for the significant

What is the nature of the problem? → Apply brain to solve → Think of the future

arm

arm

Arm Software

# Arm Forge

An interoperable toolkit for debugging and profiling

**Commercially supported by Arm**

**Fully Scalable**

**Very user-friendly**

## The de-facto standard for HPC development

- Available on the vast majority of the Top500 machines in the world
- Fully supported by Arm on x86, IBM Power, Nvidia GPUs, etc.

## State-of-the art debugging and profiling capabilities

- Powerful and in-depth error detection mechanisms (including memory debugging)
- Sampling-based profiler to identify and understand bottlenecks
- Available at any scale (from serial to parallel applications running at petascale)

## Easy to use by everyone

- Unique capabilities to simplify remote interactive sessions
- Innovative approach to present quintessential information to users

arm

# Run and ensure application correctness

Combination of debugging and re-compilation

- Ensure application correctness with **Arm DDT scalable debugger**
- Integrate with continuous integration system.
- Use version control to track changes and leverage Forge's built-in VCS support.

- Examples:
- `$> ddt --offline aprun –n 48 ./example`
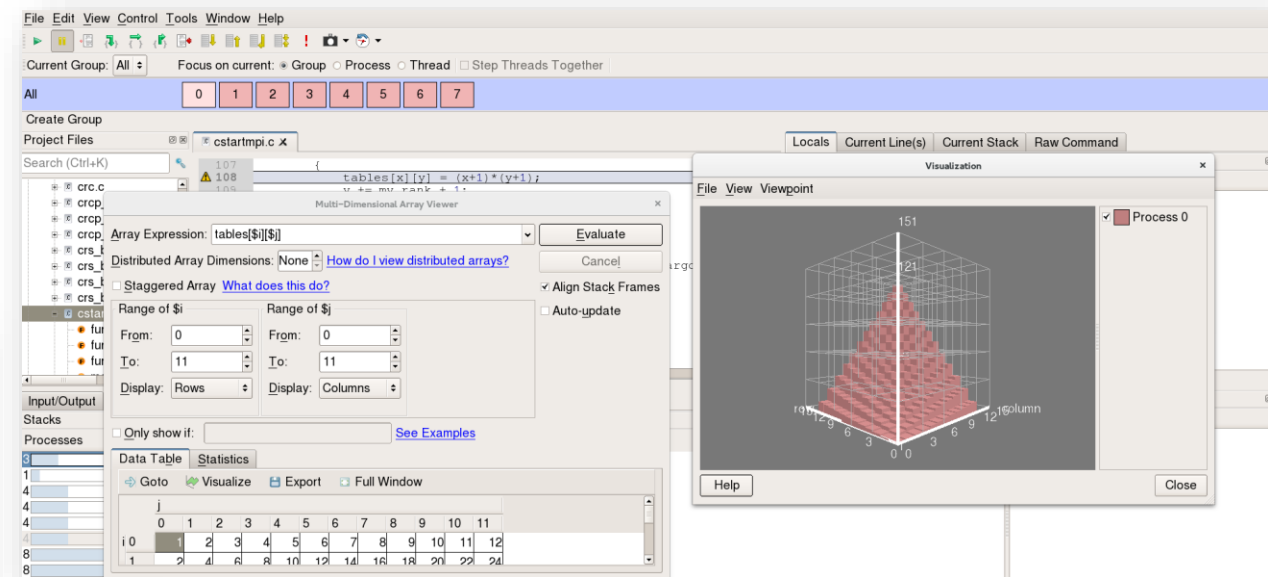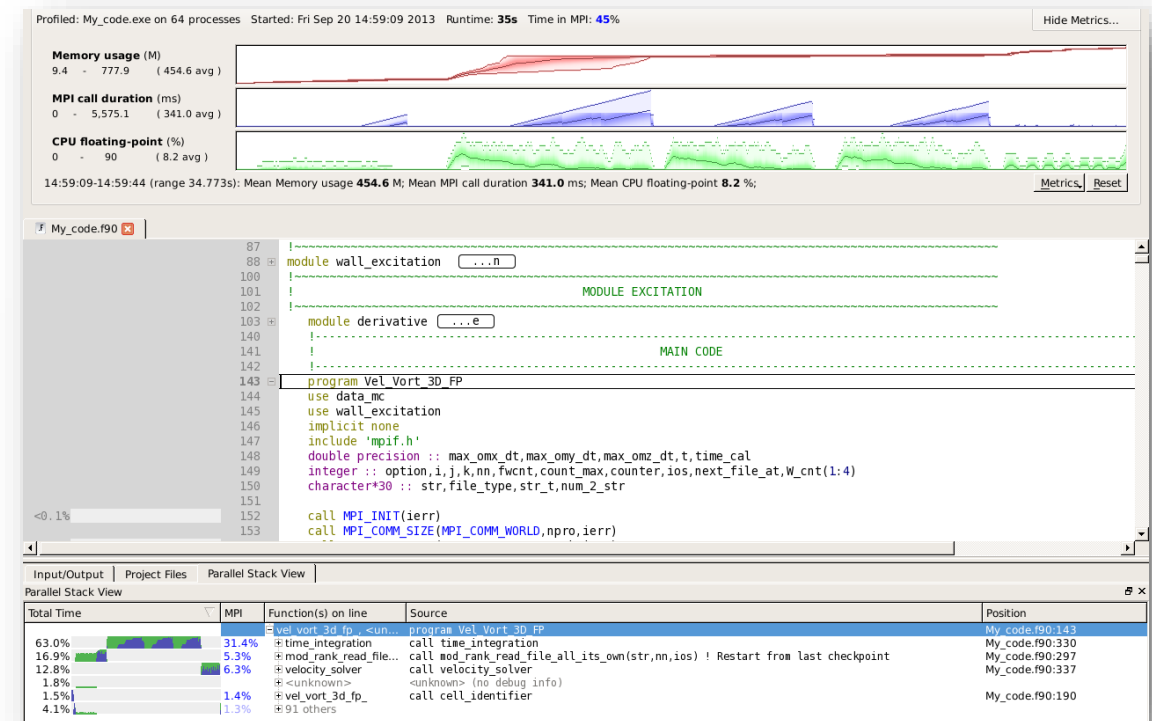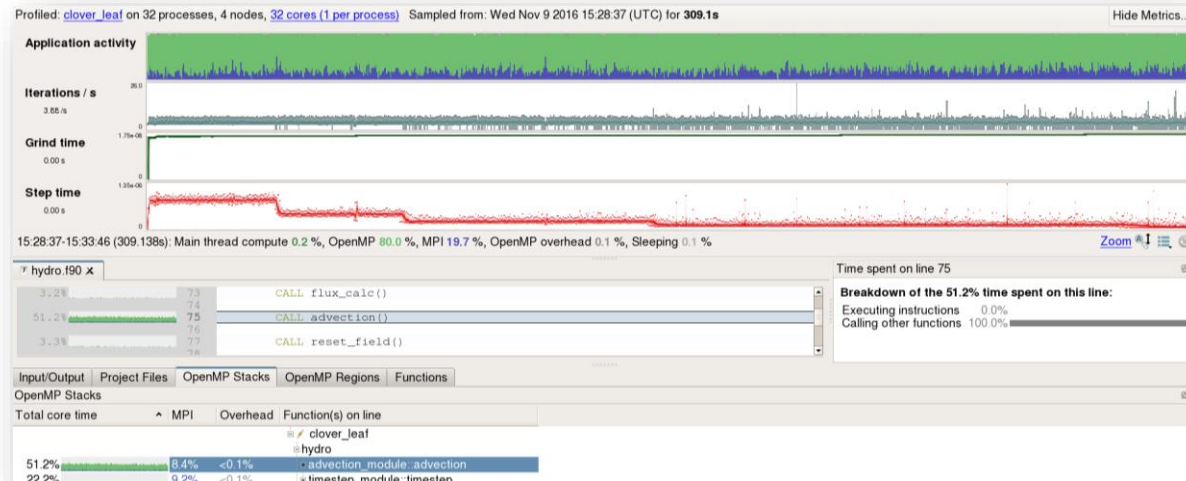- `$> ddt --connect aprun –n 48 ./example`



© 2021 A

# Visualize the performance of your application

- Measure all performance aspects with **Arm MAP parallel profiler**
- Identify bottlenecks and rewrite some code for better performance

- <u>Examples:</u>
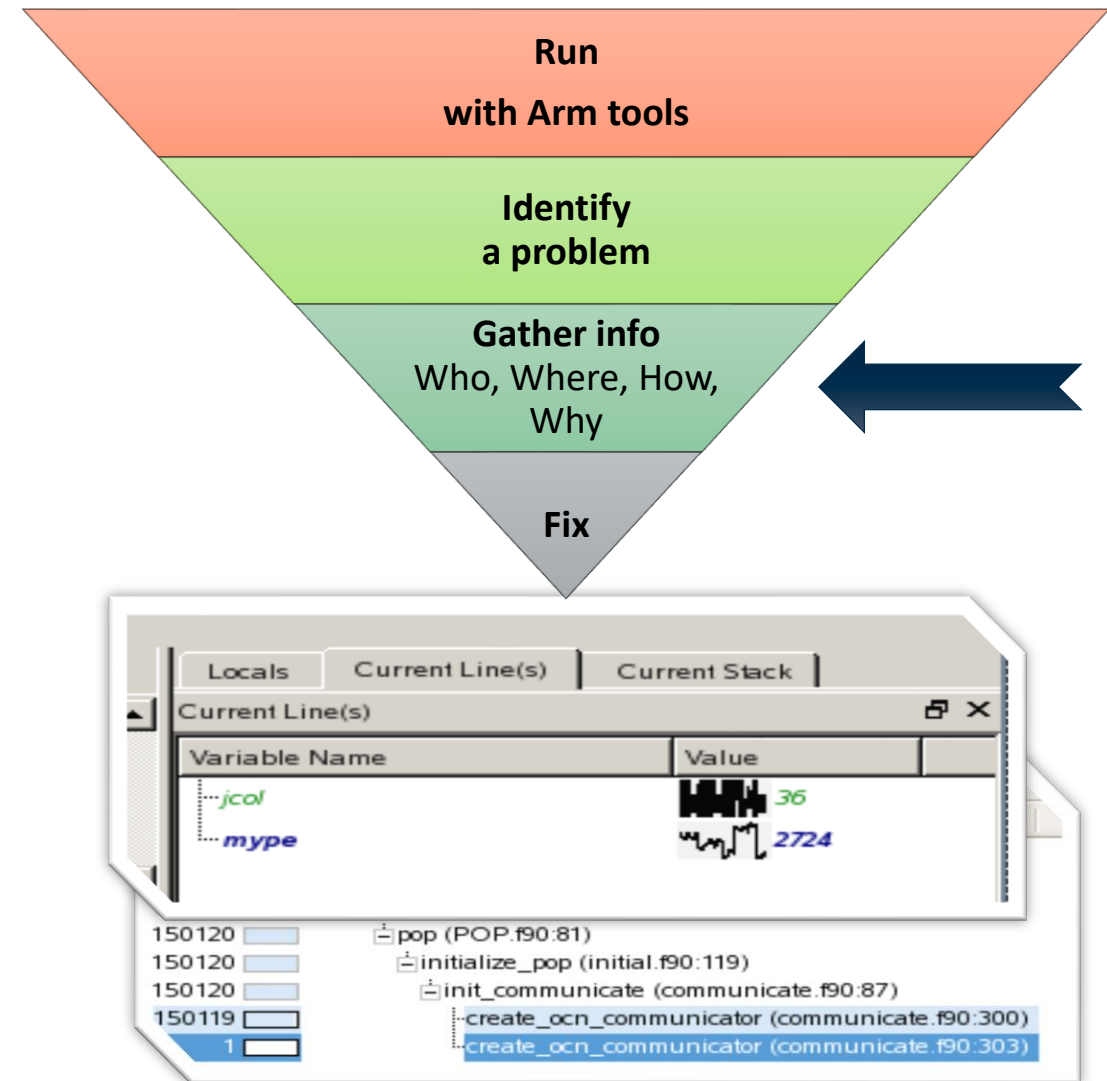- `$> map --profile -n 48 ./example`

# Debugging with DDT

# Arm DDT – The Debugger

- Who had a rogue behaviour ?
  - Merges stacks from processes and threads
- Where did it happen?
  - leaps to source
- How did it happen?
  - Diagnostic messages
  - Some faults evident instantly from source
- Why did it happen?
  - Unique "Smart Highlighting"
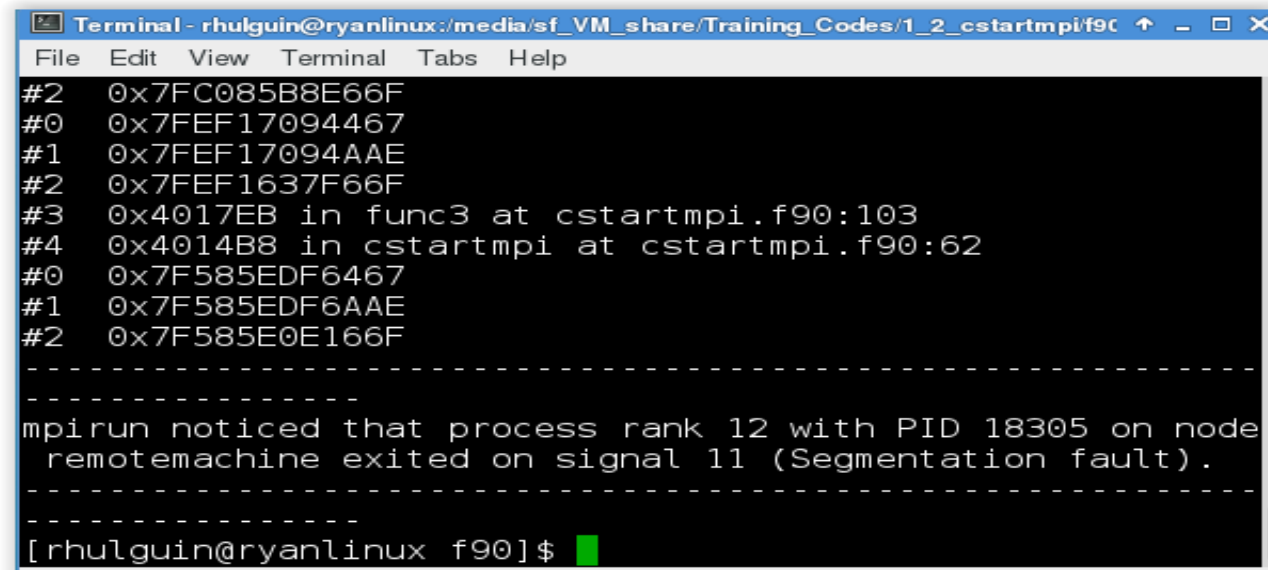  - Sparklines comparing data across processes

# Preparing Code for Use with DDT

- As with any debugger, code must be compiled with the debug flag typically **-g**

- It is recommended to turn off optimization flags i.e. **-O0**

- Leaving optimizations turned on can cause the compiler to *optimize out* some variables and even functions making it more difficult to debug
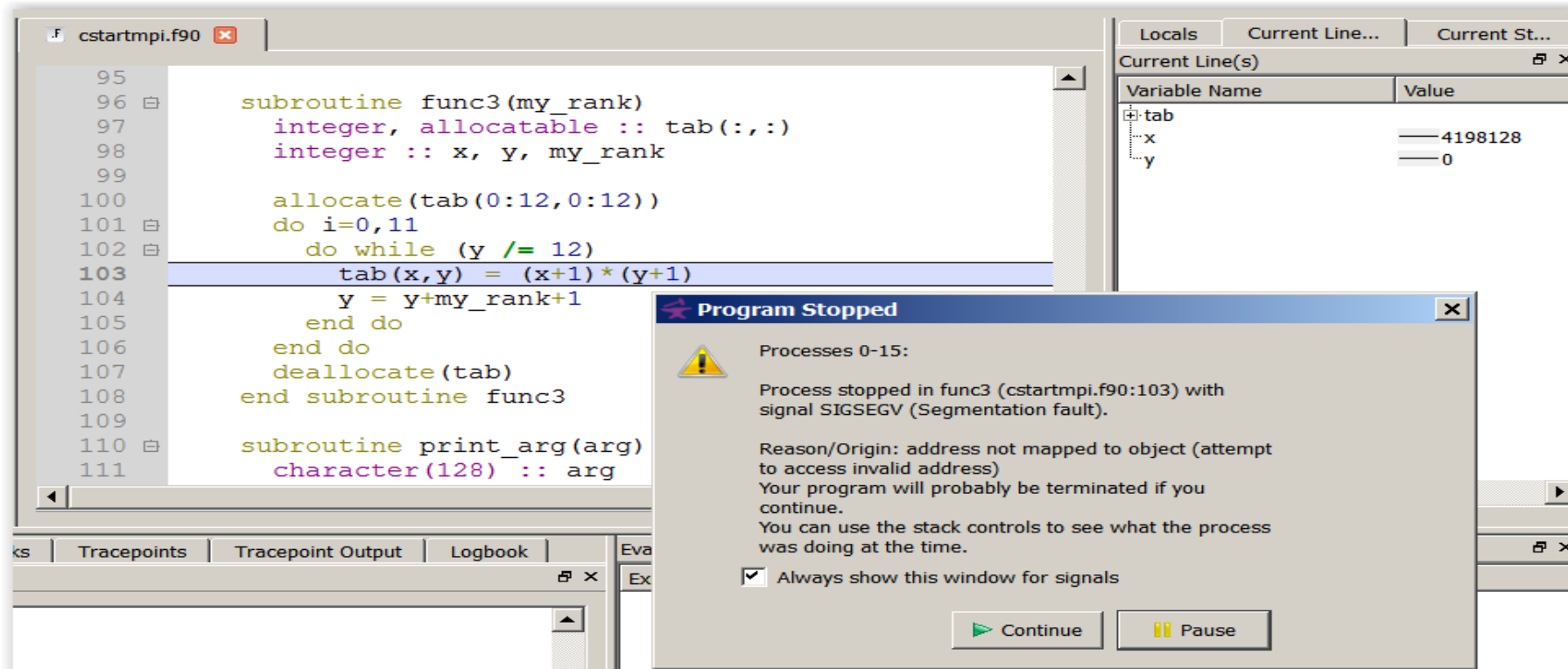
arm

# Segmentation Fault

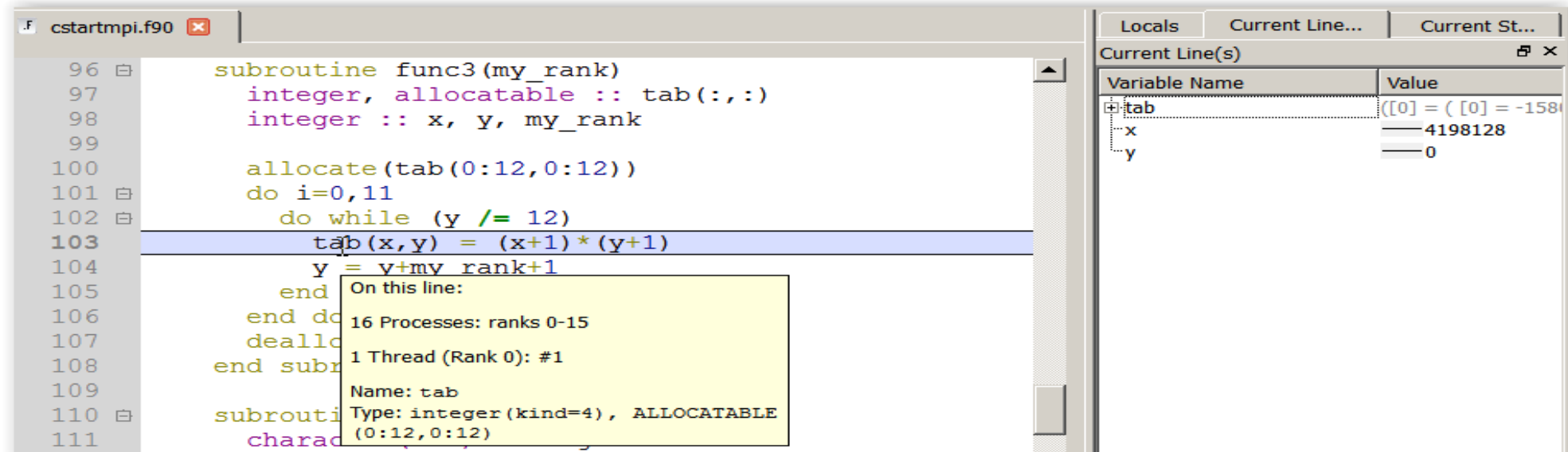- In this example, the application crashes with a segmentation error outside of DDT.



- 

- What happens when it runs under DDT?

© 2021 Arm

arm

# Segmentation Fault in DDT



- **DDT takes you to the exact line where Segmentation fault occurred, and you can pause and investigate**

# Invalid Memory Access



- The array tab is a 13x13 array, but the application is trying to write a value to tab(4198128,0) which causes the segmentation fault.
- **i** is not used, and **x** and **y** are not initialized

arm

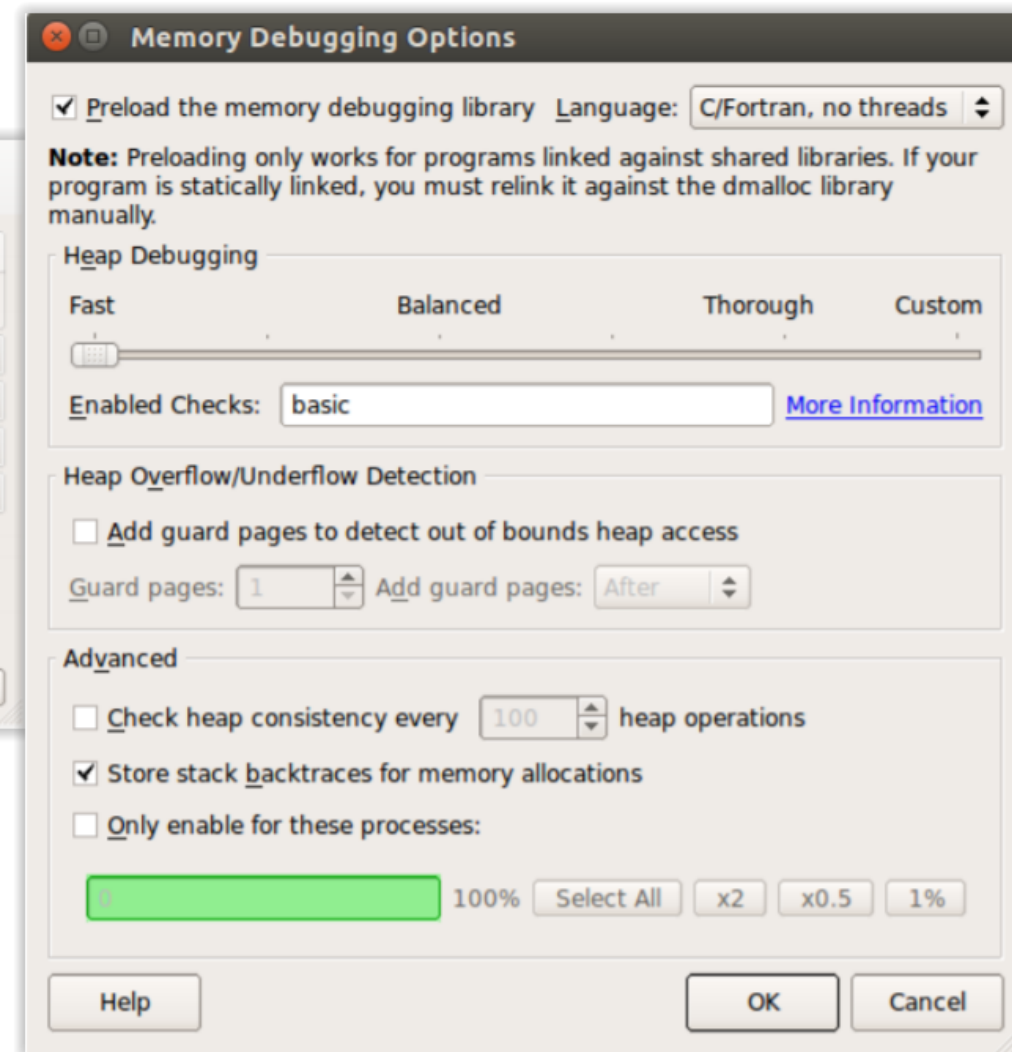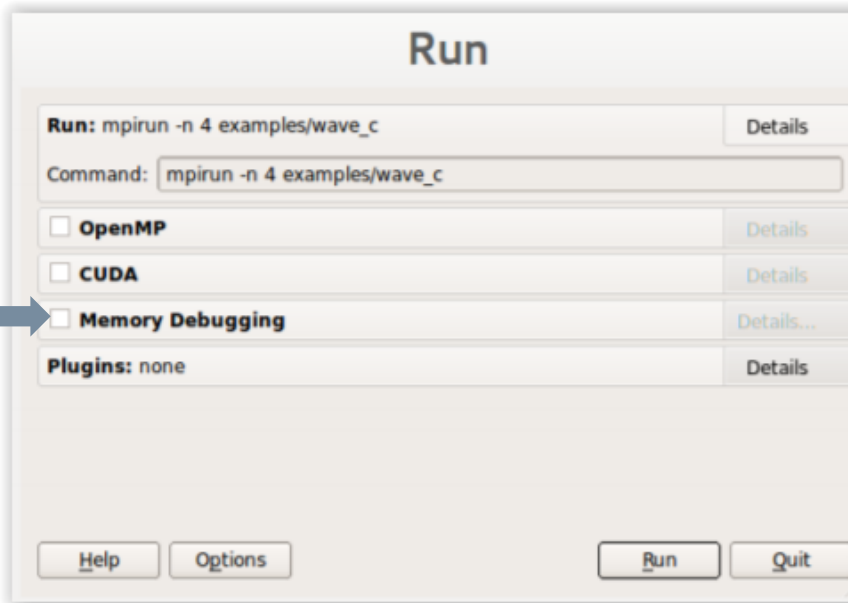# It works… Well, most of the time



**SCHRODIN BUG**

!

- A strange behaviour where the application "sometimes" crashes is a typical sign of a memory bug

- Arm DDT is able to force the crash to happen

arm

# Advanced Memory Debugging



© 2021 Arm

arm

# Heap debugging options available

**Fast**

### basic
- Detect invalid pointers passed to memory functions (e.g. malloc, free, ALLOCATE, DEALLOCATE,...)

### check-fence
- Check the end of an allocation has not been overwritten when it is freed.

### free-protect
- Protect freed memory (using hardware memory protection) so subsequent read/writes cause a fatal error.

### Added goodiness
- Memory usage, statistics, etc.

**Balanced**

### free-blank
- Overwrite the bytes of freed memory with a known value.

### alloc-blank
- Initialise the bytes of new allocations with a known value.

### check-heap
- Check for heap corruption (e.g. due to writes to invalid memory addresses).

### realloc-copy
- Always copy data to a new pointer when re-allocating a memory allocation (e.g. due to realloc)
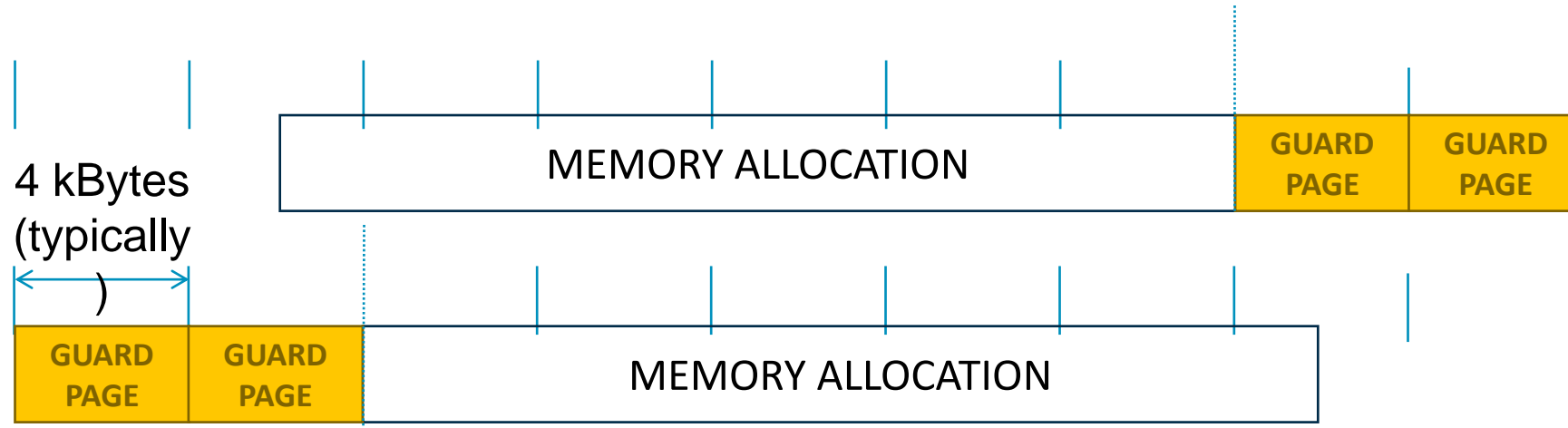
**Thorough**

### check-blank
- Check to see if space that was blanked when a pointer was allocated/freed has been overwritten.
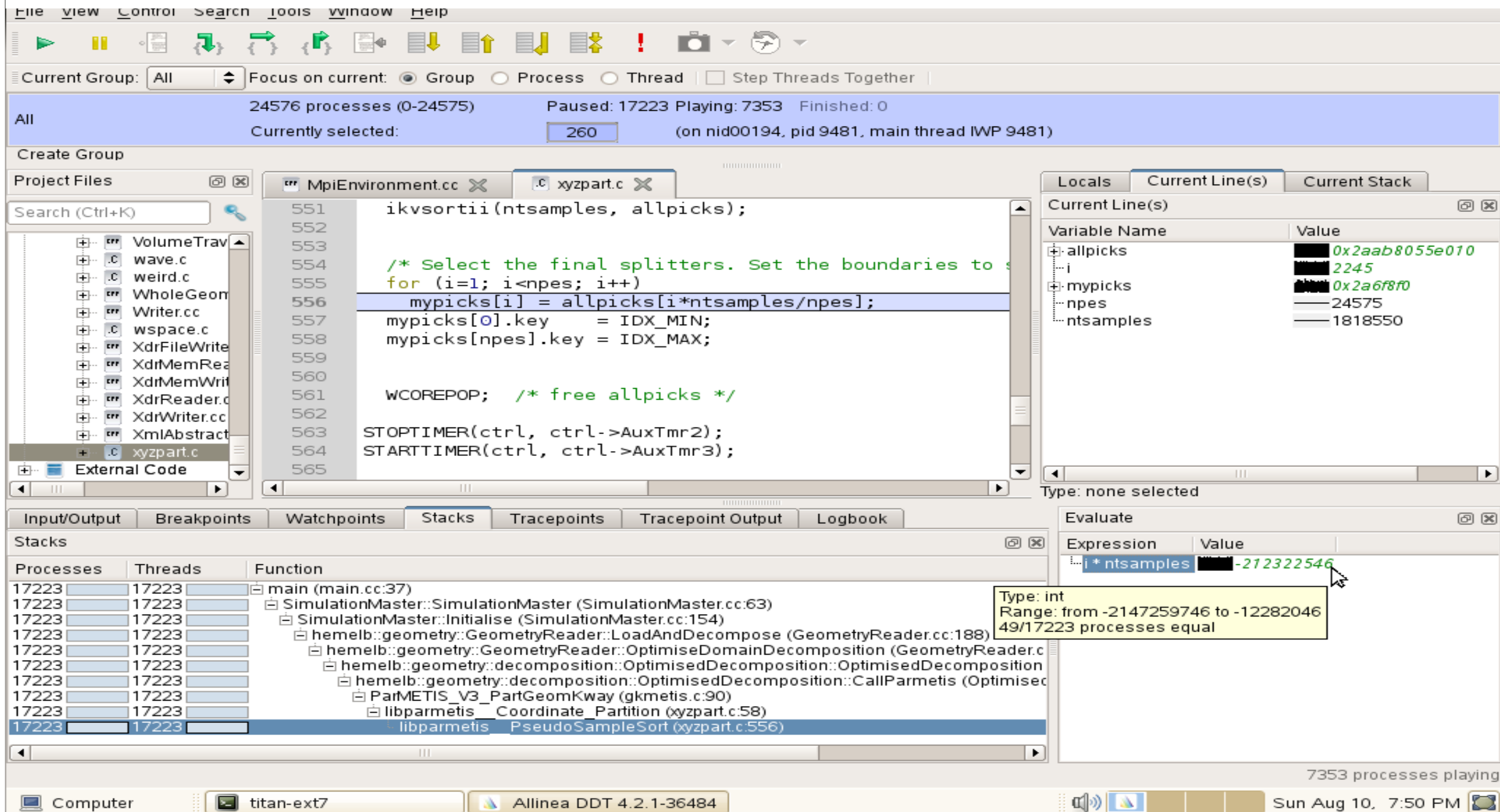
### check-funcs
- Check the arguments of addition functions (mostly string operations) for invalid pointers.
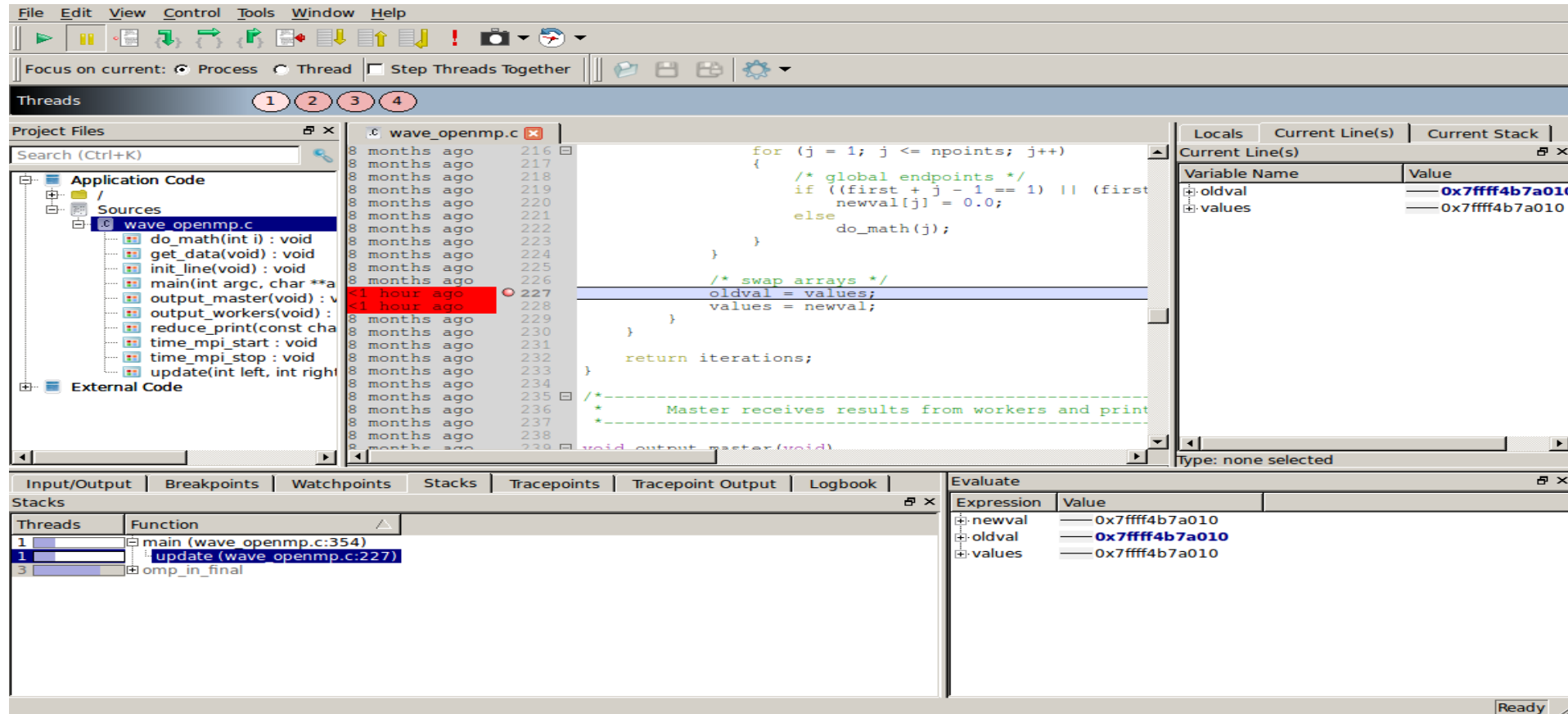
*See user-guide:*

*Chapter 12.3.2*

arm

# Guard pages (aka "Electric Fences")



4 kBytes (typically)

MEMORY ALLOCATION | GUARD PAGE | GUARD PAGE
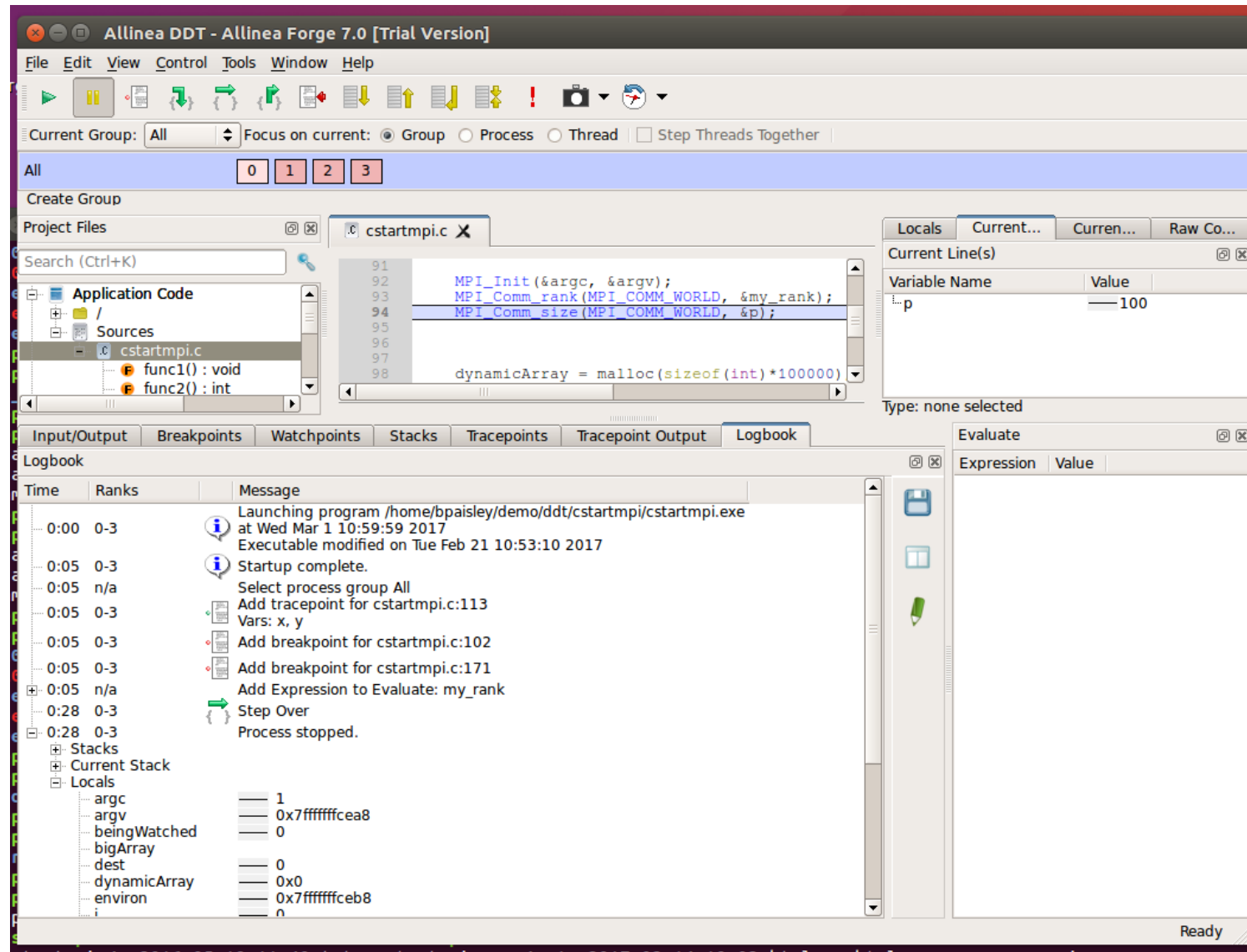
GUARD PAGE | GUARD PAGE | MEMORY ALLOCATION

- **A powerful feature…:**

    - Forbids read/write on guard pages throughout the whole execution

        *(because it overrides C Standard Memory Management library)*

- **… to be used carefully:**

    - Kernel limitation: up to 32k guard pages max ( "mprotect fails" error)

    - Beware the additional memory usage cost

arm

Current Group: All  ▼  | Focus on current: ● Group  ○ Process  ○ Thread  | ☐ Step Threads Together

**All**

24576 processes (0-24575)          Paused: 17223  Playing: 7353   Finished: 0

Currently selected:          260          (on nid00194, pid 9481, main thread IWP 9481)

Create Group

**Project Files** ▣ ✖

Search (Ctrl+K) 🔍

- ⊞ ᴱᴿᴿ VolumeTrav
- ⊞ .C wave.c
- ⊞ .C weird.c
- ⊞ ᴱᴿᴿ WholeGeom
- ⊞ ᴱᴿᴿ Writer.cc
- ⊞ .C wspace.c
- ⊞ ᴱᴿᴿ XdrFileWrite
- ⊞ ᴱᴿᴿ XdrMemRea
- ⊞ ᴱᴿᴿ XdrMemWri
- ⊞ ᴱᴿᴿ XdrReader.c
- ⊞ ᴱᴿᴿ XdrWriter.cc
- ⊞ ᴱᴿᴿ XmlAbstract
- ⊞ .C xyzpart.c
- ⊞ 🖳 External Code

**MpiEnvironment.cc** ✖  |  **xyzpart.c** ✖

```
551        ikvsortii(ntsamples, allpicks);
552
553
554        /* Select the final splitters. Set the boundaries to s
555        for (i=1; i<npes; i++)
556          mypicks[i] = allpicks[i*ntsamples/npes];
557      mypicks[0].key    = IDX_MIN;
558      mypicks[npes].key = IDX_MAX;
559
560
561      WCOREPOP;   /* free allpicks */
562
563    STOPTIMER(ctrl, ctrl->AuxTmr2);
564    STARTTIMER(ctrl, ctrl->AuxTmr3);
565
```

**Locals** | **Current Line(s)** | **Current Stack**

Current Line(s) ▣ ✖

| Variable Name | Value |
|---|---|
| ⊞ allpicks | ▓ 0x2aab8055e010 |
| ├ i | ▓ 2245 |
| ⊞ mypicks | ▓▓ 0x2a6f8f0 |
| ├ npes | ── 24575 |
| └ ntsamples | ── 1818550 |

Type: none selected

**Input/Output** | **Breakpoints** | **Watchpoints** | **Stacks** | **Tracepoints** | **Tracepoint Output** | **Logbook**

**Stacks** ▣ ✖

| Processes | Threads | Function |
|---|---|---|
| 17223 | 17223 | ⊟ main (main.cc:37) |
| 17223 | 17223 | ⊟ SimulationMaster::SimulationMaster (SimulationMaster.cc:63) |
| 17223 | 17223 | ⊟ SimulationMaster::Initialise (SimulationMaster.cc:154) |
| 17223 | 17223 | ⊟ hemelb::geometry::GeometryReader::LoadAndDecompose (GeometryReader.cc:188) |
| 17223 | 17223 | ⊟ hemelb::geometry::GeometryReader::OptimiseDomainDecomposition (GeometryReader.c |
| 17223 | 17223 | ⊟ hemelb::geometry::decomposition::OptimisedDecomposition::OptimisedDecomposition |
| 17223 | 17223 | ⊟ hemelb::geometry::decomposition::OptimisedDecomposition::CallParmetis (Optimised |
| 17223 | 17223 | ⊟ ParMETIS_V3_PartGeomKway (gkmetis.c:90) |
| 17223 | 17223 | ⊟ libparmetis__Coordinate_Partition (xyzpart.c:58) |
| 17223 | 17223 | libparmetis__PseudoSampleSort (xyzpart.c:556) |

**Evaluate** ▣ ✖

| Expression | Value |
|---|---|
| ⌐ i * ntsamples | ▓ -212322546 |

Type: int
Range: from -2147259746 to -12282046
49/17223 processes equal

7353 processes playing

# New Bugs from Latest Changes



© 2021 Arm

# Track Your Changes in a Logbook



© 2021 Arm

# Arm DDT Demo

# GPU Debugging



DDT can debug CUDA codes as well as CUDA kernels generated via OpenACC

arm

# Python Debugging



**module load intelpython36**

**module load datascience/ mpi4py-3.0.2**

ddt --connect aprun -n 8 python **%allinea_python_debug%** ./mmprod.py -s pyloop -o res -n 512

© 2021 Arm

arm

# Five great things to try with Allinea DDT

The scalable print alternative

Stop on variable change

Static analysis warnings on code errors

Detect read/write beyond array bounds

Detect stale memory allocations

arm

# Arm DDT cheat sheet

- Load the environment module
  - $ module load **forge/21.0.2**
  - $ module unload **xalt**
  - $ module unload **darshan**

- Prepare the code
  - $ cc **-O0 -g** myapp.c -o myapp.exe

- Start Arm DDT in interactive mode
  - $ **ddt** aprun -n 8 ./myapp.exe arg1 arg2

- Or use the reverse connect mechanism
  - On the login node:
    - $ ddt &
  - (or use the remote client) **<- Preferred method**
  - Then, edit the job script to run the following command and submit:
    - **ddt --connect** aprun -n 8 ./myapp.exe arg1 arg2

**arm**

# Profiling with MAP

# Arm MAP – The Profiler

**Small data files**

**<5% slowdown**

**No instrumentation**

**No recompilation**

arm

# Glean Deep Insight from our Source-Level Profiler

**Memory usage** (M)
6.1 - 62.6 ( 29.3 avg )

**MPI call duration** (ms)
0 - 727.9 ( 136.2 avg )

**MPI point-to-point** (/s)
0 - 217 ( 1.0 avg )

**MPI collectives** (/s)
0 - 172 ( 0.6 avg )

**CPU memory access** (%)
0 - 100 ( 18.7 avg )

**CPU floating-point** (%)
0 - 100 ( 10.9 avg )

**CPU vector** (%)
0 - 100 ( 37.6 avg )

**CPU branch** (%)
0 - 60 ( 0.5 avg )

Track memory usage across the entire application over time

Spot MPI and OpenMP imbalance and overhead

Optimize CPU memory and vectorization in loops

Detect and diagnose I/O bottlenecks at real scale

arm

# Profile of 2d Laplace Solver with Jacobi Iteration



© 2021 Arm

arm

# Tracking Largest Change

- `// Compare newly computed value with old value`
- `diff = fabs(grid2[i][j] – grid1[i][j]);`
- `// Track largest change between new and old values`
- `maxDiff = diff > maxDiff ? Diff : maxDiff;`

<br>

- `If (diff > maxDiff)`
- `    then maxDiff= diff;`
- `Else`
- `   maxDiff = maxDiff;`

© 2021 Arm

**arm**

# Conditional Removal from Innermost Loop



**20 % faster, also operation is now vectorized**

© 2021 Arm

arm

# Initial profile of CloverLeaf shows surprisingly unequal I/O

Each I/O operation should take about the same time, but it's not the case.



© 2021 Arm

# Symptoms and causes of the I/O issues
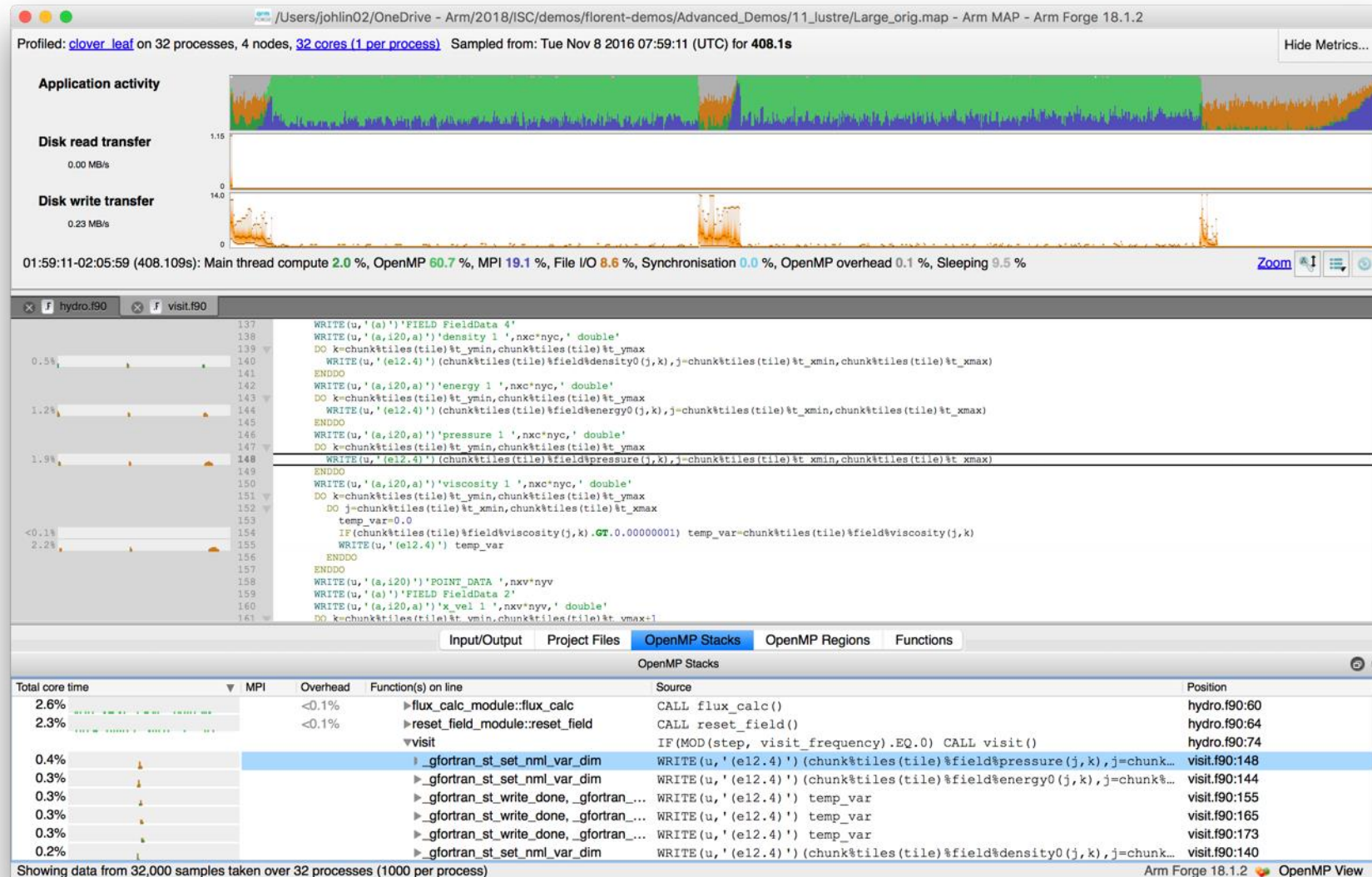
Sub-optimal file format and surprise buffering.



- Write rate is less than 14MB/s.

- Writing an ASCII output file.

- Writes not being flushed until buffer is full.

  - Some ranks have much less buffered data than others.

  - Ranks with small buffers wait in barrier for other ranks to finish flushing their buffers.

arm

# Solution: use HDF5 to write binary files

Using a library optimized for HPC I/O improves performance and portability.



© 2021 Arm

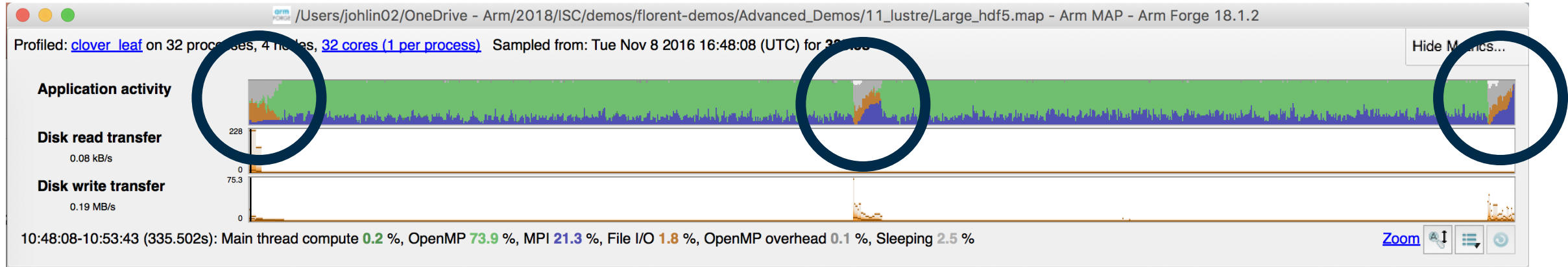# Solution: use HDF5 to write binary files

Using a library optimized for HPC I/O improves performance and portability.



- Replace Fortran write statements with HDF5 library calls.

  - Binary format reduces write volume and can improve data precision.

  - Maximum transfer rate now 75.3 MB/s, over 5x faster.

- Note MPI costs (blue) in the I/O region, so room for improvement.
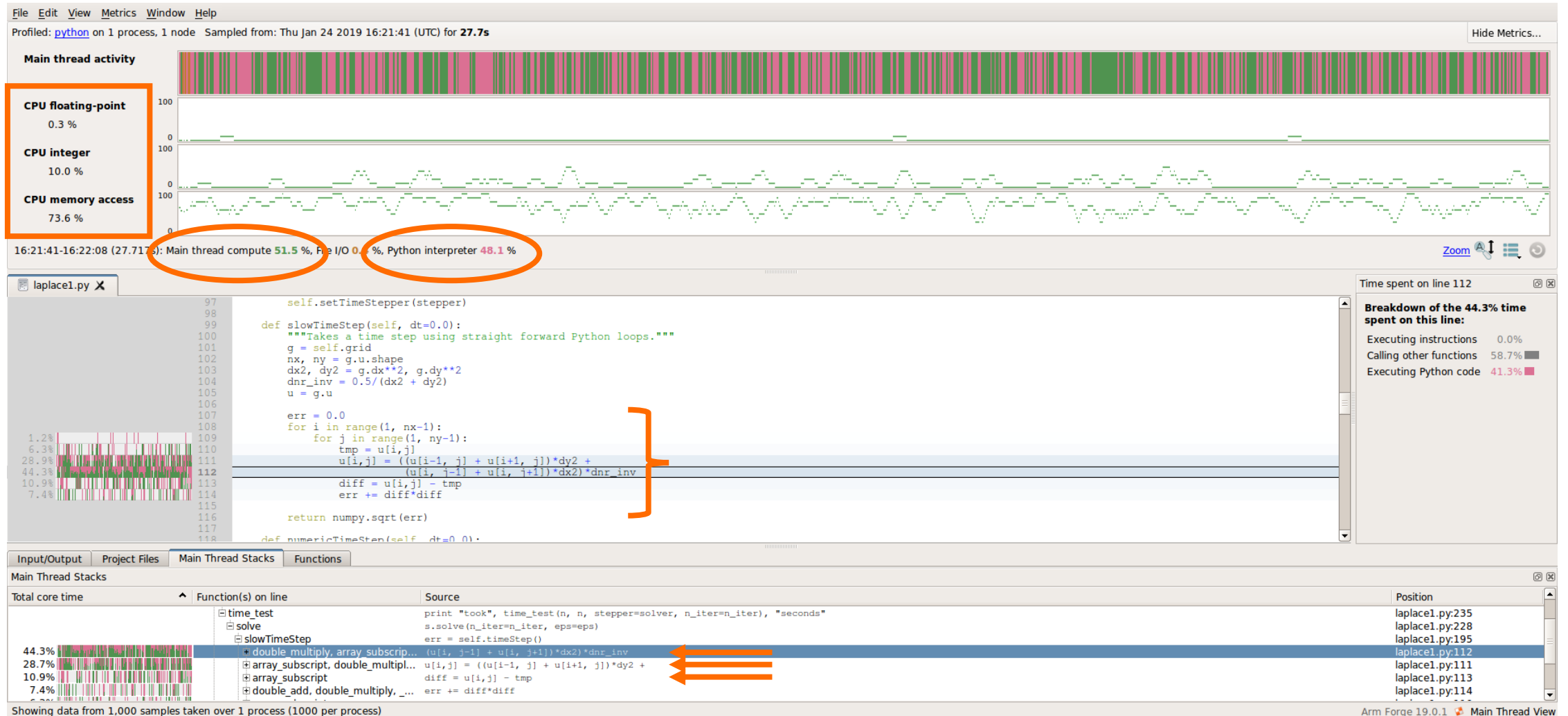
arm

# Arm MAP: Python profiling

- Launch command
  - $ **python** ./laplace1.py slow 100 100

- Profiling command
  - $ **map --profile python** ./laplace1.py slow 100 100
  - --profile: non-interactive mode
  - --output: name of output file

- Display profiling results
  - $ **map** laplace1.map

**Laplace1.py**

```python
[…]
err = 0.0
for i in range(1, nx-1):
    for j in range(1, ny-1):
        tmp = u[i,j]
        u[i,j] = ((u[i-1, j] + u[i+1, j])*dy2 +
            (u[i, j-1] + u[i, j+1])*dx2)*dnr_inv
        diff = u[i,j] - tmp
        err += diff*diff
return numpy.sqrt(err)
[…]
```

arm

# Naïve Python loop (laplace1.py slow 100 1000)



© 2021 Arm

arm

# Optimizing computation on NumPy arrays

## Naïve Python loop

```python
err = 0.0
for i in range(1, nx-1):
  for j in range(1, ny-1):
    tmp = u[i,j]
    u[i,j] = ((u[i-1, j] + u[i+1, j])*dy2 +
    (u[i, j-1] + u[i, j+1])*dx2)*dnr_inv
    diff = u[i,j] - tmp
    err += diff*diff
return numpy.sqrt(err)
```
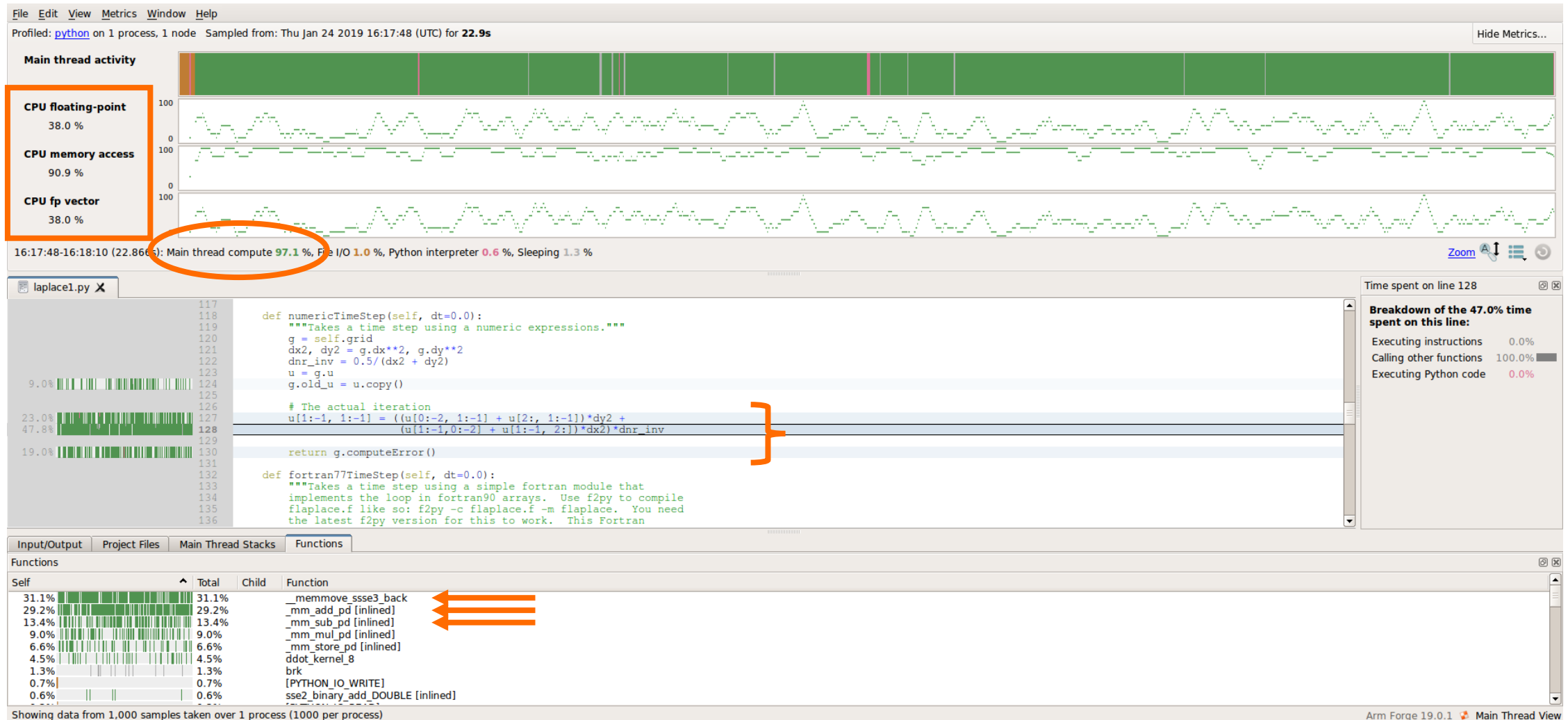
## NumPy loop

```python
u[1:-1, 1:-1] =
    ((u[0:-2, 1:-1] + u[2:, 1:-1])*dy2 +
    (u[1:-1,0:-2] + u[1:-1, 2:])*dx2)*dnr_inv

return g.computeError()
```

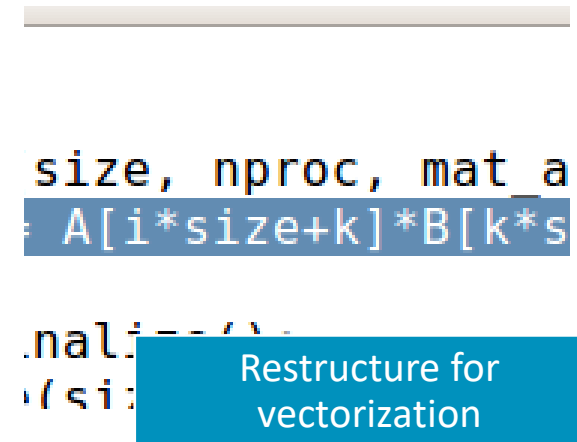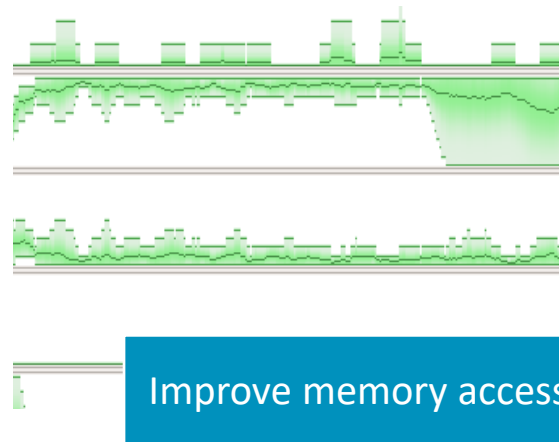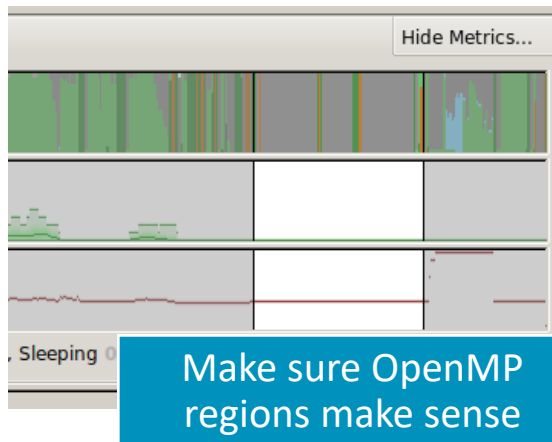**arm**

# NumPy array notation (laplace1.py numeric 1000 1000)

**This is 10 times more iterations than was computed in the previous profile**



© 2021 Arm

# Arm MAP cheat sheet

- Load the environment module (manually specify version)
  - $ module load **forge/21.0.2**

- Unload Darshan module (It wraps MPI calls which cannot be used with MAP)
  - $ module unload darshan

- Compile the code with dynamic linking
  - $ cc -O3 **-g -dynamic** myapp.c -o myapp.exe
  - Edit the job script to run Arm MAP in "profile" mode
  - $ **map --profile** aprun -n 8 ./myapp.exe arg1 arg2

- Open the results
  - On the login node:
    - $ map myapp_Xp_Yn_YYYY-MM-DD_HH-MM.map
  - (or load the corresponding file using the remote client connected to the remote system or locally)

**arm**

# Six Great Things to Try with Allinea MAP



Find the peak memory use

Fix an MPI imbalance

Remove I/O bottleneck

Make sure OpenMP regions make sense

Improve memory access

Restructure for vectorization

arm

# Theta Specific Settings

# Configure the remote client

**Install the Arm Remote Client**

- Go to : https://developer.arm.com/products/software-development-tools/hpc/downloads/download-arm-forge

**Connect to the cluster with the remote client**

- Open your Remote Client
- Create a new connection: Remote Launch ➜ Configure ➜ Add
  - Hostname: `<username>@theta.alcf.anl.gov`
  - Remote installation directory:
    `/soft/debuggers/soft/debuggers/forge-21.0.2-2021-06-11`

**Create a separate connection if using thetagpu**

- Create a new connection: Remote Launch ➜ Configure ➜ Add
  - Hostname: `<username>@theta.alcf.anl.gov thetagpusn1`
  - Remote installation directory:
    `/soft/debuggers/soft/debuggers/forge-21.0.2-2021-06-11`

arm

# Static Linking Extra Steps

- To enable advanced memory debugging features in DDT with static binaries, you must link explicitly against our memory debugging libraries

- Simply add the link flags to your Makefile, or however appropriate

- 

  lflags = -L/soft/debuggers/ddt/lib/64 -Wl,--undefined=malloc -ldmallocthcxx -Wl,--allow-multiple-definition

arm

# Debugging on Thetagpu

- The latest Forge modules are not available on thetagpu, but you can you use the installed software directly

- Use the temporary license shown below
  ```
  export ALLINEA_FORCE_LICENCE_FILE=/grand/ATPESC2021/EXAMPLES/track-6-tools/arm_forge/Licence.trial
  ```

- Debug your GPU code using
  ```
  /lus/theta-fs0/software/debuggers/forge-21.0.2-2021-06-11/bin/ddt --connect gpu_code.exe
  ```

**arm**

# Profiling on Theta

- Although static binaries are created by default on Theta, it is recommended to build dynamic executables for profiling purposes with the compiler flag **-dynamic**

- If you get library missing errors, reload the intel module

- **module unload intel**

- **module load intel**


- If you get GdbmiParser errors set the following environment variable

- **export ALLINEA_FORCE_DEBUGGER=gdb-82**

© 2021 Arm

arm

# Examples for hands-on session

- Examples are available at `/grand/ATPESC2021/EXAMPLES/track-6-tools/arm_forge`

- README files are available for each example

**arm**

# arm

# Questions?

# arm

Thank You
Danke
Gracias
谢谢
ありがとう
Asante
Merci
감사합니다
धन्यवाद
Kiitos
شكرًا
ধন্যবাদ
תודה

# arm